

Introdução à Ciência de Dados

Cayan Portela

UniCEUB

March 3, 2023

Definição

Amazon Web Services (AWS)

"Uma abordagem multidisciplinar que combina princípios e práticas das áreas de matemática, estatística, inteligência artificial e engenharia da computação para analisar grandes quantidades de dados."

- Analisar dados
- Detecção de padrões
- Extrair informação

Onde estão os dados?



- Diversas fontes de dados
 - → Operações financeiras.
 - \rightarrow e-commerce, web data.
 - → Repositórios de documentos.
 - → Redes sociais.

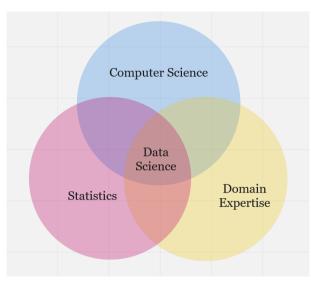
- Diversas fontes de dados
 - → Operações financeiras.
 - → e-commerce, web data.
 - → Repositórios de documentos.
 - → Redes sociais.
- Como estão sendo utilizados?
 - → Otimização de recursos.
 - → Recomendação.
 - → Tomadas de decisão.
 - → Personalização.

■ Exemplos de casos de uso

- → Aprovar ou não um empréstimo.
- → Ofertas de produtos.
- → Engajamento de usuários.
- \rightarrow Teste AB.
- → Jurimetria.

Ciência de Dados





- Matemática e Estatística
 - → Cálculo e Álgebra Linear (como otimizar uma função?).
 - → Probabilidade e Estatística.
 - → Inferência e Modelagem preditiva.



- Matemática e Estatística
 - → Cálculo e Álgebra Linear (como otimizar uma função?).
 - → Probabilidade e Estatística.
 - → Inferência e Modelagem preditiva.
- Computação
 - → Cloud (GCP, AWS, Azure, etc)
 - → Linguagens de programação: Python, R, SQL.
 - → Ferramentas de versionamento (trabalho em grupo)
 - → Pipelines e steps de produtização.



Matemática e Estatística

- → Cálculo e Álgebra Linear (como otimizar uma função?).
- → Probabilidade e Estatística.
- → Inferência e Modelagem preditiva.

■ Computação

- → Cloud (GCP, AWS, Azure, etc)
- → Linguagens de programação: Python, R, SQL.
- → Ferramentas de versionamento (trabalho em grupo)
- → Pipelines e steps de produtização.

■ Expertise na Área

- → Qual métrica avaliar?
- → Qual alteração realizar?
- → Trade-off entre complexidade e resultados.
- → Entendimento de todo o ecossistema do produto/serviço.

Matemática e Estatística



- Cálculo e Álgebra Linear
 - → Operações matriciais.
 - → Derivada e Integral.
 - → Ponto ótimo em superfície de respostas.

- Atividades de um cientista de dados
 - → Cloud (GCP, AWS, Azure, etc)
 - → Linguagens de programação: Python, R, SQL.
 - → Versionamento: git.
 - → IDE: VScode, RStudio, PyCharm, spyder...
 - → Etapas de modelagem.







DATA SCIENCE



MACHINE LEARNING ENGINEERING



DATA ENGINEERING



what it is

tools/languages

Data scientists use statistics to build models that help companies draw insights and make predictions from their data.

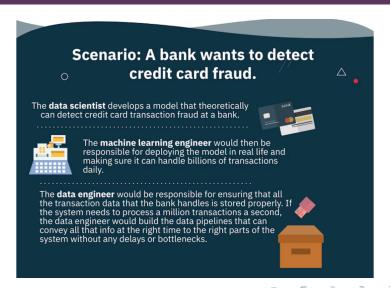
Programming languages like **Python** and **R**. Data science libraries like **pandas**, **scikit-learn** and **iupyter notebooks**. Machine Learning Engineering (MLE) is the art and science of deploying models developed by data scientists and turning them into a live production system.

Tools for model implementation like **TensorFlow**. Tools for model deployment like **Microsoft** Azure, Amazon SageMaker, Google Cloud ML. Data engineers set up the infrastructure for others to work on; they are responsible for data storage, data transportation. etc.

Data storage and pipeline tools like **Oracle**, NoSQL tools like **Cassandra**, queuing and messaging systems like **Kafka**, and workflow tools like **Airflow**









UniCEUB



UniCEUB

Alavancar o crescimento com estratégias baseadas em dados

Aquisição



UniCEUB

- Aquisição
- Engajamento



UniCEUB

- Aquisição
- Engajamento
- Retenção



UniCEUB

- Aquisição
- Engajamento
- Retenção
- Crescimento

Alavancar o crescimento com estratégias baseadas em dados

- Aquisição
- Engajamento
- Retenção
- Crescimento

→ Crescimento da base ativa.



UniCEUB



UniCEUB

- Facebook
 - \rightarrow 7 amigos em 10 dias.

- Facebook
 - \rightarrow 7 amigos em 10 dias.
- Twitter
 - → Seguir 30 pessoas.

- Facebook
 - \rightarrow 7 amigos em 10 dias.
- Twitter
 - → Seguir 30 pessoas.
- Slack
 - → Trocar 2000 mensagens

- Facebook
 - \rightarrow 7 amigos em 10 dias.
- Twitter
 - → Seguir 30 pessoas.
- Slack
 - → Trocar 2000 mensagens
- Netflix
 - → Assistir algo em 90 segundos

Melhoria contínua



UniCEUB

Ciência de dados viabiliza adaptação e evolução.

Melhoria contínua



UniCEUR

Ciência de dados viabiliza adaptação e evolução.

- Netflix
 - → Assistir algo em 90 segundos

Como melhorar o processo?

Melhoria contínua



UniCEUR

Ciência de dados viabiliza adaptação e evolução.

- Netflix
 - → Assistir algo em 90 segundos

Como melhorar o processo?

- → Experimentação.
- → Conclusão a partir de dados.





Cells	Cell 1 (Control)	Cell 2	Cell 3
Box Art	A NETELIX ORIGINAL SHORT GAME	SHORT GAME	SHORT
	Default artwork	14% better take rate	6% better take rate





Ciência de dados



UniCEUR

Transformar dados em informação

- Analisar dados para obter insights
- Identificar padrões, tendencias, correlacoes.
- Entender e contextualizar os dados.

Ciência de dados



UniCEUB

Variável aleatória

Etimologicamente, Variável aleatória é uma:

Ciência de dados



UniCEUB

Variável aleatória

Etimologicamente, Variável aleatória é uma:

■ Variável: Possui um valor desconhecido

Variável aleatória

Etimologicamente, Variável aleatória é uma:

- Variável: Possui um valor desconhecido
- Aleatória: Pode possuir valores diferentes, com diferentes probabilidades

Em casos aplicados, nossas "variáveis aleatórias" são as informações observadas.

Tipos de dados



- Quantitativos
 - → Discreta
 - → Contínua
- Qualitativos
 - → Nominal
 - \rightarrow Ordinal

UniCEUR

Dados quantitativos assumem valores numéricos em sua forma natural (dados numéricos).

- Discretos
 - Assume valores inteiros. Números naturais
 - Ex: Número de "caras" após 3 lançamentos de uma moeda.
- Contínuos
 - Pode assumir qualquer valor no espaço amostral.
 - Ex: Quantidade de água em uma garrafa.

Dados qualitativos assumem valores categóricos em sua forma natural (categorias, classes)

Não possui significado matemático

- Nominal
 - Não possui hierarquia.
 - Gênero, Nacionalidade.
- Ordinal
 - Possui relação entre valores assumidos.
 - Grau de instrução.
 - Em geral, podemos ter variáveis latentes classificadas como dados ordinais (Ex: avaliação de satisfação)

Estruturados x Não Estruturados

Estruturados

- Informação organizada
- armazenamento mais leve
 - Tabelas (csv, txt, excel)
 - Bancos de dados relacionais

Não Estruturados

- Estrutura diversificada
- armazenamento mais pesado
 - **Imagens**
 - Videos
 - Texto

Como os dados estao gravados? Arquivos e extensões

- .csv (comma-separated-values)
 - → ', '; '
- .txt (text)
 - → tab, espaço
- excel
 - → .xls, .xlsx
 - ightarrow tambem suportam csv
- json
 - → chave valor

Dados Tabulares



UniCEUB

Forma mais comum de dados estruturados

Tabela: linhas e colunas

Em geral:

■ Linhas: observações (ex: cada pessoa)

■ Colunas: variaveis (ex: atributos de cada pessoa)

Exemplo

■ Considere a seguinte tabela:

Ramo	Vendas (R\$)	Aluguel (R\$)	Nº aluguéis	Tipo
Doces	2.000	300	24	Quiosque
Papelaria	3.000	700	6	Loja
Magica	1.500	600	2	Quiosque
Amendoim	5.000	1000	36	Loja
Pipoca	4.000	800	30	Loja

■ Considere a seguinte tabela:

Ramo	Vendas (R\$)	Aluguel (R\$)	Nº aluguéis	Tipo
Doces	2.000	300	24	Quiosque
Papelaria	3.000	700	6	Loja
Magica	1.500	600	2	Quiosque
Amendoim	5.000	1000	36	Loja
Pipoca	4.000	800	30	Loja

Classifique as variáveis da tabela.



■ Considere a seguinte tabela:

Ramo	Vendas (R\$)	Aluguel (R\$)	Nº aluguéis	Tipo
Doces	2.000	300	24	Quiosque
Papelaria	3.000	700	6	Loja
Magica	1.500	600	2	Quiosque
Amendoim	5.000	1000	36	Loja
Pipoca	4.000	800	30	Loja

Classifique as variáveis da tabela.

O que são parecem ser as observações?

Preço (R\$)	Tipo Quarto	Quarto Dividido	Capacidade
194	Quarto Privativo	Não	2 pessoas
264	Quarto Privativo	Não	2 pessoas
344	Casa toda	Não	4 pessoas
433	Quarto Privativo	Não	2 pessoas
485	Casa Toda	Não	6 pessoas

■ Considere a seguinte tabela:

Preço (R\$)	Tipo Quarto	Quarto Dividido	Capacidade
194	Quarto Privativo	Não	2 pessoas
264	Quarto Privativo	Não	2 pessoas
344	Casa toda	Não	4 pessoas
433	Quarto Privativo	Não	2 pessoas
485	Casa Toda	Não	6 pessoas

Classifique as variáveis da tabela.

O que são as observações?

Dados Tabulares



UniCEUB

Importar e manipular

- R
- python
- excel

Devemos acessar os dados para manipular e trabalhar

Medidas resumo



UniCEUB

Não vamos trabalhar com um único valor

22

Não vamos trabalhar com um único valor

22

E sim com vários valores

22, 22, 23, 24, 25, 25, 26, 26, 27, 39, 59, 79

Não vamos trabalhar com um único valor

22

E sim com vários valores

22, 22, 23, 24, 25, 25, 26, 26, 27, 39, 59, 79

Suponha que a sequencia acima represente uma amostra aleatória contendo a idade de 12 pessoas. Que informações podemos extrair?

→ Média

$$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

- → Mediana
 - Valor central
- → Moda
 - Valor com maior frequência
- Amplitude
 - Diferença entre o maior valor e o menor valor.

Considerando a idade das 12 pessoas na amostra coletada, calcule: média, moda, mediana, amplitude.

22, 22, 23, 24, 25, 25, 26, 26, 27, 39, 59, 79

Exercícios



UniCEUR

Com uma amostra de tamanho 12, é simples calcular na mão. Para lidar com uma grande quantidade de dados, precisamos utilizar ferramentas disponíveis: R, python, VSCode, RStudio, SQL...